

Automated Generation of Control Concepts Annotation Rules Using Inductive Logic Programming

April 24, 2022

Basel Shbita University of Southern California / GE Research (intern) Abha Moitra GE Research

As part of the Graphical Symbol Expression Network (GraSEN) project DARPA's Recovery of Symbolic Mathematics from Code (ReMath) program

Intro

- Cyber-Physical Systems (CPS)
 - Physical systems in which mechanism is controlled by computer-based algorithms
 - Include micro-controller software
 - Control theory connects applied math to CPS engineering
 - Primitive and higher-level (composite) mathematical concepts
 - Constant, gain, saturation
 - Output signal, error signal
 - PID controller





Automated Generation of Control Concepts Annotation Rules Using Inductive Logic Programmin

Intro



Intro

- Cyber-Physical Systems (CPS)
 - Physical systems in which mechanism is controlled by computer-based algorithms
 - Include micro-controller software
 - Control theory connects applied math to CPS engineering
 - Primitive and higher-level (composite) mathematical concepts
- Bridging the gap between binary code analysts & domain SMEs
- Binary \rightarrow code analysis framework \rightarrow Expert
- How?
 - binary decompilation
 - semantic code pattern analysis





Automated Generation of Control Concepts Annotation Rules Using Inductive Logic Programmin

Active

ARM platform

binary

Expert

x86 platform

binary

Goal

Semantic Model



knowledge acquisition & discovery to assist the SME?



Inductive Logic Programming (ILP)

• ILP

- Advantages:
 - Requires small amounts of data
 - Uses expressive FOL (First-Order Logic)
 - Takes background knowledge into account
- Used for classification & prediction
 - Bioinformatics, NLP
- Given background knowledge (B) & specific observations (E+, E-) induce general rules (hypothesis, H)
- Aleph, an ILP system
- A Learning Engine for Proposing Hypotheses
- Learning from entailments (LFE)







Problem Definition

- Generate classification (annotation) rules for math & control concepts using ILP
 - Input
 - (OWL) data representing basic code elements in a CPS program (e.g., variables, functions, operators)
 - Positive & negative examples the objective concept (identified by the SME)
 - Output
 - (Prolog) program describing the objective concept (e.g., constants, error signals, PI controller)
 - Iterative process if the generated rule identifies false positive \rightarrow add example & repeat





Approach





Generating the ILP Data (owl2aleph)





Generating the ILP Data (owl2aleph)



Generating the ILP Data (owl2aleph)

			Clear Generate Exam	generate .f & .n files	A			
Po	sitives:			Negatives:				
newfeature(xf3407a4281654856a4f1d newfeature(xf899593b6b954b098b7a!	1692e51d5 9fdf8aba0b	8f9). oab).	newf	newfeature(xb1b9c84b70824e668b52a527a20f9298).				
positive hyperedges			ne,	egative hyperedges				
HyperEdge	line start	ftype	lambda	add as positive	d to			
HyperEdge xf849fbef3f59455c8d9b0aae61a7fe02	line start 16	ftype xassign	lambda lambda timestep,error, integrator_state)	add as positive , integrator_state: ((timestep * error) +	d to			
HyperEdge xf849fbef3f59455c8d9b0aae61a7fe02 xfbcab097a142468380f366753096577e	line start 16 18	ftype xassign xassign	<pre>lambda lambda timestep,error, integrator_state) lambda error,Kp_M,Ki_N integrator_state))</pre>	add as positive add as positive f, integrator_state: ((timestep * error) + M, integrator_state: ((error * Kp_M) + (Ki_M *	d to			



Rule Generation from ILP Data (SWI-Prolog)

Bottom (most-specific) clause





Evaluation

• Dataset

- three OWL files
- 8,974 triples
- 61 math & control instances
- 9 different classes of concepts

• Results

- 7/9: perfect F1 score
 - Small training data
 - Significant reduction in # literals
- 2/9: low F1
 - not enough data to separate positives from negatives
- Time complexity is excellent < 1s

[
$\operatorname{Concept}$	(E+,	Bottom	Reduced	Time	True	False	False	Precision	Recall	F1
	E-)	clause	clause	[sec-	posi-	posi-	nega-			
		size	size	onds]	tives	tives	tives			
Difference	(2, 1)	24	3	0.063	4	0	0	1.0	1.0	1.0
Sum	(2, 2)	29	3	0.063	9	0	0	1.0	1.0	1.0
Gain	(2, 1)	23	3	0.047	9	0	0	1.0	1.0	1.0
Division	(2, 2)	24	3	0.094	3	0	0	1.0	1.0	1.0
Constant	(2, 5)	11	6	0.125	23	0	0	1.0	1.0	1.0
Error signal	(2, 5)	24	6	0.859	2	0	0	1.0	1.0	1.0
PI controller	(2, 5)	45	5	0.453	3	0	0	1.0	1.0	1.0
Output signal	(2, 3)	12	12	0.859	3	3	0	0.50	1.0	0.67
Reference signal	(2, 2)	11	11	0.375	3	17	0	0.15	1.0	0.26
Switch										
Magnitude	Not enough positives $(E+)$									
saturation					0	(• /			
PID controller										



Evaluation											
LVAIGATION	Concept	(E+, E-)	Bottom clause size	Reduced clause size	Time [sec- onds]	True posi- tives	False posi- tives	False nega- tives	Precision	Recall	F1
Dataset	Difference	(2.1)	24	3	0.063	4	0	0	1.0	1.0	1.0
• three OW gain(A) :-					063	9	0	0	1.0	1.0	1.0
• 8 97/ trip xfunction(A,	B), has ope	erato	r mult	(B).	047	9	0	0	1.0	1.0	1.0
 9 different classes of concepts Desults 	outputs(A xfunctior	• A,B), N(A,C)	var_e), fun	kplici c_not_	t(B), in_lo	, var_ pop_b]	_assig lock(C	ned_o	nce(B),		1.0 1.0 1.0
Results	Output signal	(2, 3)	12	12	0.859	3	3	0	0.50	1.0	0.67
 7/9: perfect F1 score Small training data Significant reduction in # literals 	Reference signal Switch Magnitude saturation	a1 $(2, 2)$ 11 11 0.375 3 17 0 0.15 1.0									0.26
• 2/9: low F1		-									



Discussion

- Platform is feasible & effective in terms of time, completeness & robustness
 - Results are **satisfying** (simple & composite concepts)
 - Quality of ILP generated rules depends on supplied input
- Still, many challenges exist:
 - Local vs. global variables (structs, enumerations, etc...)
 - In-line code (implicit vs. explicit)
 - Pointers
 - Implementation "style" (e.g., saturation: if/else vs. arithmetic)
- Looking forward:
 - Rules can be **folded** back into the semantic model
 - Learn multiple level of knowledge
 - Use probabilistic soft-logic on instance data to introduce confidence
 - Joint-classification
 - Incorporate higher level (or system) knowledge (e.g., there is a single integrator in the code)
 - Augmentation
 - Layout features
 - Memory allocations



Conclusion

- An automated platform for generating classification rules for math & control concepts
 - SME not required to know language formalism to describe knowledge in detail
 - Fast, iterative induction with human-in-the-loop
 - Structured & semantic
 - Easily extended & refined, automatically or manually
 - Handy during design/engineering
 - Incorporate feedback
 - Generate explanations



